# Limits and Convergence properties of the

# Sequentially Markovian Coalescent

Thibaut Sellinger[1*], Diala Abu Awad[1], Aurélien Tellier[1]

[1] Professorship for Population Genetics,

Department of Life Science Systems,

Technical University of Munich

[*] Corresponding author, thibaut.sellinger@tum.de

**1**       **Abstract**

**2**    Many methods based on the Sequentially Markovian Coalescent (SMC)

**3**    have been and are being developed. These methods make use of genome

**4**    sequence data to uncover population demographic history. More recently,

**5**    new methods have extended the original theoretical framework, allowing

**6**    the simultaneous estimation of the demographic history and other biolog-

**7**    ical variables. These methods can be applied to many different species,

**8**    under different model assumptions, in hopes of unlocking the popula-

**9**    tion/species evolutionary history. Although convergence proofs in par-

**10**    ticular cases have been given using simulated data, a clear outline of the

**11**    performance limits of these methods is lacking. We here explore the limits

**12**    of this methodology, as well as present a tool that can be used to help

**13**    users quantify what information can be confidently retrieved from given

**14**    datasets. In addition, we study the consequences for inference accuracy

**15**    violating the hypotheses and the assumptions of SMC approaches, such

**16**    as the presence of transposable elements, variable recombination and mu-

**17**    tation rates along the sequence and SNP call errors. We also provide a

**18**    new interpretation of the SMC through the use of the estimated transi-

**19**    tion matrix and offer recommendations for the most efficient use of these

**20**    methods under budget constraints, notably through the building of data

**21**    sets that would be better adapted for the biological question at hand.

**22**    ***Keywords***— Hidden Markov Model, Ancestral Recombination Graph, Popu-

**23**    lation Genetics

# 1  Introduction

Recovering the demographic history of a population has become a central theme in evolutionary biology. The demographic history (the variation of effective population size over time) is linked to environmental and demographic changes that existing and/or extinct species have experienced (population expansion, colonization of new habitats, past bottlenecks) [14, 43, 4]. Current statistical tools to estimate the demographic history rely on genomic data [51] and these inferences are often linked to archaeological or climatic data, providing novel insights on the evolutionary history [70, 33, 44, 1, 12, 26, 25]. From these analyses, evidence for migration events have been uncovered [26, 5], as have genomic consequences of human activities on other species [9]. Linking demographic history to climate and environmental data greatly supports the field of conservation genetics [10, 17, 40]. Indeed, using such approaches can help ecologists in detecting effective population size decrease [68], and thus serve as a guide in maintaining or avoiding the erosion of genetic diversity in endangered populations, and potentially predicting the consequences of climate change on genetic diversity [27]. In addition, studying the demographic histories of different species in relation to one another can unveil latent biological or environmental evolutionary forces [16], unveiling links and changes within entire ecosystems. With the increased accuracy of current methods, the availability of very large and diverse data sets and the development of new theoretical frameworks, the demographic history has become an information that is essential in the field of evolution [48, 6]. However, obtaining unbiased estimations/interpretations of the demographic history remain challenging [3, 8].

The most sophisticated methods to infer demographic history make use of

2

48 whole genome polymorphism data. Among the state-of-the-art methods, are those

49 based on the theory of the Sequentially Markovian Coalescent (SMC) developed by

50 McVean and Cardin[35] after the work of Wiuf and Hein [69], corrected by Marjoram

51 and Wall [31] and first applied to whole genome sequences by Li and Durbin [26], who

52 introduced the now well-known, Pairwise Sequentially Markovian Coalescent (PSMC)

53 method. PSMC allows demographic inference of populations with unprecedented ac-

54 curacy, while requiring only one sequenced diploid individual. This method uses the

55 distribution of SNPs along the genome between the two sequences to account for and

56 infer recombination and demographic history of a given population, assuming neu-

57 trality and panmixia. Although PSMC was a breakthrough in demographic inference,

58 it has limited power in inferring more recent events. In order to address this issue,

59 PSMC has been extended to account for multiple sequences (*i.e.* more than two) into

60 the method known as the Multiple Sequentially Markovian Coalescent (MSMC) [50].

61 By using more sequences, MSMC better infers recent events and also provides the pos-

62 sibility of inferring population splits using the cross-coalescent rate. MSMC, unlike

63 PSMC, is not based on SMC theory [35] but on SMC' theory [31], therefore MSMC

64 applied to only two sequences has been defined as PSMC'. Methods developed after

65 MSMC followed suit, with MSMC2 [30] extending PSMC by incorporating pairwise

66 analysis, increasing efficiency and the number of sequences that can be inputted (up to

67 a hundred), resulting in more accurate results. SMC++ [63] brings the SMC theory to

68 another level by allowing the use of hundreds of unphased sequences (MSMC requires

69 phased input data) and breaking the piece-wise constant population size hypothesis,

70 while accounting for the sample frequency spectrum (SFS). Because SMC++ incorpo-

71 rates the SFS in the estimation of demographic history, it increases accuracy in recent

72 time [63]. SMC++ is currently the state of the art SMC-based method for big data

73  sets (>20 sequences), but seems to be outperformed by PSMC when using smaller

74  data sets [45]. In a similar vein, the Ascertained Sequentially Markovian Coalescent

75  (ASMC) [42] extends the SMC theory to estimate coalescence times at the locus scale

76  from ascertained SNP array data, something that was made possible by the theory

77  presented by Hobolth and Jensen [18].

78      More recently, a second generation of SMC-based methods have been devel-

79  oped. New features have been added to the initial SMC theory, extending its ap-

80  plication beyond simply inferring past demography [1, 53, 66]. The development of

81  C-PSMC [16] allows the interpretation of estimated demographic history in the light

82  of coevolution between species, making the first link between demographic history es-

83  timated by PSMC and evolutionary forces (although biological interpretation remains

84  limited). iSMC [1] extends the PSMC theory to account and infer the variation of the

85  recombination rate along sequences, unlocking recombination map estimations. An im-

86  pressive advancement is the development of MSMC-IM, which to some extent solves

87  the population structure problem, allowing the accurate and simultaneous inference of

88  the demographic history and population admixture [66]. eSMC [53] incorporates com-

89  mon biological traits (such as self-fertilization and dormancy) and demonstrated the

90  strong effect life-history traits can have on demographic history estimations. Results

91  which could not be explained under the initial SMC hypotheses can now be explained

92  by the potential presence of measurable phenomena.

93      New methods have been developed since PSMC, that have been either strongly

94  inspired by the SMC [54, 62] or that are completely dissociated from it [58, 2, 49, 20,

95  29, 19, 57, 65]. Though there are alternative approaches, methods based on the SMC

96  are still considered state of the art, and remain widely used [32, 3, 59], notably in

4

97 human evolution studies [59, 45]. However, each described method has its specificity,

98 being based on different hypothesis in order to solve a particular problem or shortcom-

99 ings of existing methodology. Although all these methods allow a new and different

100 interpretation of genomic data, none of these methods guarantees unbiased inference,

101 and their limitations have underlined how crucial and challenging demographic infer-

102 ence is, highlighting the complementarity and usefulness of applying several inference

103 methods on a given dataset.

104     SMC-based methods display very good fits when using simulated data, espe-

105 cially when using simple single population models based on typical human data param-

106 eters [63, 50, 53, 66]. However, the SMC makes a large number of hypotheses [26, 50]

107 that are often violated in data obtained from natural populations. When inputting

108 data from natural populations, extracting information or correctly interpreting the

109 results can become troublesome [8, 64, 3] and several studies address the consequences

110 of hypothesis violation [15, 8, 49, 34, 52]. They bring to light how strongly population

111 structure or introgression influence demographic history estimation if not correctly ac-

112 counted for [15, 8]. Furthermore, some SMC-based methods require phased data (such

113 as MSMC [50] and MSMC-IM [66]), and phasing errors can lead to a strong overes-

114 timation of population size in recent time [63]. The effect of sequencing coverage has

115 also been tested in Nadachowska et al. [37], showing the importance of high coverage

116 in order to obtain trustworthy results, and yet, SMC methods seem robust to genome

117 quality [45]. Selection, if not accounted for, can result in a bottleneck signature [52],

118 and there is currently no solution to this issue within the SMC theory, though it could

119 be addressed using different theoretical frameworks that are being developed [55, 38].

120 More problematic, is the ratio of effective recombination over effective mutation rates

**121** $\frac{\rho}{\theta}$, which, if it is greater than one, biases estimations [63, 1, 53]. It is also important to

**122** keep in mind that there can be deviations between $\frac{\rho}{\theta}$ and the ratio of recombination

**123** rate over mutation rate measured experimentally ($\frac{r}{\mu}$), as the former can be greatly

**124** influenced by life-history, such as in organisms displaying self-fertilization, partheno-

**125** genesis or dormancy, and this can lead to issues when interpreting results (*e.g.* [53]).

**126** It is thus necessary to keep in mind that the accuracy of SMC-based methods depends

**127** on which of the many underlying hypothesis are prone to being violated by the data

**128** sets being used.

**129**     In an attempt to complement previous works, we here study the limits and

**130** convergence properties of methods based on the Sequentially Markovian Coalescent.

**131** We first define the limits of SMC-based methods (*i.e.* how well they perform theo-

**132** retically), which we will call the best-case convergence. In order to do this, we use

**133** a similar approach to [13, 41, 19], and compare simulation results obtained with the

**134** simulated Ancestral Recombination Graph (ARG) as input to results obtained from

**135** sequences simulated under the same ARG, so as to study the convergence properties

**136** linked to data sets in the absence of hypothesis violation. We test several scenarios

**137** to check whether there are instances, where even without violating the underlying hy-

**138** potheses of the methodology, the demographic scenarios cannot be retrieved because

**139** of theoretical limits (and not issues linked with data). We also study the effect of the

**140** optimization function (or composite likelihood) and the time window of the analysis

**141** on the estimations of different variables. Lastly, we test the effect of commonly vi-

**142** olated hypotheses, such as the effect of the variation of recombination and mutation

**143** rates along the sequence and between scaffolds, errors in SNP calls and the presence

**144** of transposable elements and link abnormal results to specific hypothesis violations.

6

145 Through this work, our aim is to provide guidelines concerning the interpretation of

146 results when applying this methodology on data sets that may violate the underlying

147 hypotheses of the SMC framework.

# 148  2   Materials and Methods

149 In this study we use four different SMC-based methods: MSMC, MSMC2, SMC++

150 and eSMC. All methods are Hidden Markov Models and use whole genome sequence

151 polymorphism data. The hidden states of these methods are the coalescence times

152 (or genealogies) of the sample. In order to have a finite number of hidden states,

153 they are grouped into $x$ bins ($x$ being the number of hidden states). The reasons for

154 our model choices are as follows: $i$) MSMC, unlike any other method, focuses on the

155 first coalescence event of a sample of size $n$, and thus exhibits different convergence

156 properties [50], $ii$) MSMC2 computes coalescence times of all pairwise analysis from

157 a sample of size $n$, and can deal with a large range of data sets [58], $iii$) SMC++

158 [63] is the most advanced and efficient SMC method and lastly, $iv$) eSMC [53] is a re-

159 implementation of PSMC' (similar to MSMC2), which will contribute to highlighting

160 the importance of algorithmic translations as it is very flexible in its use and outputs

161 intermediate results necessary for this study.

## 162  2.1   SMC methods

### 163  2.1.1   PSMC', MSMC2 and eSMC

164 PSMC' and methods that stem from it (such as MSMC2 [30] and eSMC [53]) focus on

165 the coalescence events between only two individuals (or sequences in practice), and,

7

166  as a result, do not require phased data. The algorithm goes along the sequence and

167  estimates the coalescence time at each position. In order to do this, it checks whether

168  the two sequences are similar or different at each position. The presence or absence of

169  a segregating site along the sequence (determined by the population mutation rate $\theta$)

170  is used to infer the hidden state (*i.e.* coalescence time). However, the hidden state is

171  only allowed to change in the event of a recombination, which leads to a break in the

172  current genealogy. Thus, the population recombination rate $\rho$ constrains the inferred

173  changes of hidden states along the sequence (for a detailed description of the algorithm

174  see [50, 66, 53]).

### 2.1.2  MSMC

176  MSMC is mathematically and conceptually very similar to the PSMC' method. Unlike

177  other SMC methods, it simultaneously analyses multiple sequences and because of

178  this, MSMC requires the data to be phased. In combination with a second HMM,

179  to estimate the external branch length of the genealogy, it can follow the distribution

180  of the first coalescence event in the sample along the sequences. However, due to

181  computational load, MSMC cannot analyze more than 10 sequences simultaneously

182  (for a detailed description see [50]).

### 2.1.3  SMC++

184  SMC++ is slightly more complex than MSMC or PSMC. Though it is conceptually

185  very similar to PSMC', mathematically it is quite different. SMC++ has a differ-

186  ent emission matrix compared to previous methods because it calculates the sample

187  frequency spectrum of sample size $n + 2$, conditioned on the coalescence time of two

188  "distinguished" haploids and $n$ "undistinguished" haploids. In addition SMC++ offers

189  features such as a cubic spline to estimate demographic history (*i.e.* not a piece-wise

190  constant population size). The SMC++ algorithm is fully described in [63].

### 191  2.1.4   Best-case convergence

192  Using sequence simulators such as msprime [21] or scrm [60], one can simulate the

193  Ancestral Recombination Graph (ARG) of a sample. Usually the ARG is given through

194  a sequence of genealogies (*e.g.* a sequence of trees in Newick format). From this ARG,

195  one can find what state of the HMM the sample is in at each position. Hence, one

196  can build the series of states along the genomes, and build the transition matrix.

197  The transition matrix, is a square matrix of dimension $x$ (where $x$ is the number

198  of hidden states) counting all the possible pairwise transitions between the $x$ states

199  (including from a given state to itself). Using the transition matrix built directly

200  from the exact ARG, one can estimate parameters using eSMC or MSMC as if they

201  could correctly infer the hidden states. Hence estimations using the exact transition

202  matrix represents the upper bound of performance for these methods. We choose

203  to call this upper bound the best-case convergence (since it can never be reached in

204  practice). For this study's purpose, a second version of the R package eSMC [53]

205  was developed. This package enables the building of the transition matrix (for eSMC

206  or MSMC), and can then use it to infer the demographic history. The package is

207  mathematically identical to the previous version, but includes extra functions, features

208  and new outputs necessary for this study. The package and its description can be found

209  at https://github.com/TPPSellinger/eSMC2.

## 2.1.5    Baum-Welch algorithm

211    SMC-based methods can use different optimization functions to infer the demographic

212    parameters ( *i.e.* likelihood or composite likelihood).  The four studied methods use

213    the Baum-Welch algorithm to maximize the likelihood.  MSMC2 and SMC++ imple-

214    ment the original Baum-Welch algorithm (which we call the complete Baum-Welch

215    algorithm), whereas eSMC and MSMC compute the expected composite likelihood

216    $Q(\theta|\theta^t)$ based only on the transition matrix (which we call the incomplete Baum-

217    Welch algorithm).  The use of the complete Baum-Welch algorithm or the incomplete

218    one can be specified in the eSMC package.  The composite likelihood for SMC++ and

219    MSMC2 is given by equations 1 and the composite likelihood for eSMC and MSMC

220    by equation 2:

$$Q(\Theta|\Theta^t) = \nu_{\Theta^t} log(P(X_1|\Theta)) + \sum_{X,Y} E(X,Z|\Theta^t) log(P(X|Z,\Theta)) + \sum_{X,Y} E(Y,X|\Theta^t) log(P(Y|X,\Theta))$$

(1)

221        and :

$$Q(\Theta|\Theta^t) = \sum_{X,Y} E(X,Z|\Theta^t) log(P(X|Z,\Theta)),$$  (2)

222    with:

223        • $\nu_\Theta$ : The equilibrium probability conditional to the set of parameters $\Theta$.

224        • $P(X_1|\Theta)$ :  The probability of the first hidden state conditional to the set of

225            parameters $\Theta$.

226        • $E(X,Z|\Theta^t)$ :  The expected number of transitions of X from Z conditional to

227            the observation and set of parameters $\Theta^t$.

10

228  • $P(X|Z, \Theta)$ : The transition probability from state Z to state X, conditional to

229      the set of parameters $\Theta$.

230  • $E(Y, X|\Theta^t)$ The expected number of observations of type Y that occurred during

231      state X conditional to observation and set of parameters $\Theta^t$.

232  • $P(Y|X, \Theta)$ : The emission probability conditional to the set of parameters $\Theta$.

### 2.1.6  Time window

234  Each tested SMC-based method has its own specific time window for which estima-

235  tions are made. Note that hidden states are defined as discretized intervals, as a

236  consequences of which the boundaries, length and number of states of the time win-

237  dow do implicitly affect inferences. For example, the original PSMC has a time window

238  wider than PSMC', so that estimations cannot be compared one to one. To measure

239  the effect of choosing different time window parameters, we analyze the same data

240  with four different settings. The first time window is the one used for PSMC' defined

241  in [50]. The second time window is that of MSMC2 [66] (similar to the one of the

242  original PSMC [26]), which we call "big" since it goes further in the past and in more

243  recent time than that of PSMC'. We then define a time window equivalent to the first

244  one (i.e. PSMC') shifted by a factor five in the past (first time window, *i.e.* hidden

245  states, multiplied by five). The last one is a time window equivalent to the first one

246  (i.e. PSMC') shifted by a factor five in recent time (i.e. first time window divided by

247  five).

### 2.1.7  Regularization penalty

249  To avoid inferring unrealistic demographic histories with variations of population size

250  that are too strong and/or too rapid, SMC++ introduced a regularization penalty

11

251 (https://github.com/popgenmethods/smcpp). This parameter penalizes population

252 size variation. In SMC++, the lower value of the penalty the more the estimated size

253 history is a line (*i.e.* constant population size). Regularization penalty was also imple-

254 mented in eSMC. Setting the regularization penalty parameter to 0 is equivalent to no

255 penalization, and the higher the parameter value, the more population size variations

256 are penalized (https://github.com/TPPSellinger/eSMC2 for more details). We tested

257 the effect of regularization on inferences with both methods using simulated sequence

258 data. The sequence data was simulated under sawtooth demographic histories with

259 different amplitudes of population size variation.

260     All the command lines to analyze the data generated can be found in S2 of the

261 Appendix.

## 2.2 Simulated sequence data

263 Throughout this paper we simulate different demographic scenarios using either the

264 coalescence simulation program scrm [60] or msprime [21]. We use scrm for the best-

265 case convergence as it can output the genealogies in a Newick format (which we use as

266 input). We use scrm, which outputs simulated sequences in the ms format, to simulate

267 data for eSMC, MSMC, MSMC2. We use msprime to simulate data for SMC++ since

268 msprime is more efficient than scrm for big sample sizes [21] and can directly output

269 .vcf files (which is the input format of SMC++).

### 2.2.1 Absence of hypothesis violation

271 We simulate five different demographic scenarios: saw-tooth (successions of population

272 size exponential expansion and decrease), bottleneck, exponential expansion , expo-

273 nential decrease and constant population size. Each of the scenarios with varying

12

274 population size is tested under four amplitude parameters (*i.e.* by how many fold the

275 population size varies: 2, 5, 10, 50). We infer the best-case convergence under four

276 different sequence lengths ($10^7$, $10^8$, $10^9$ and $10^{10}$ bp) and choose the per site mutation

277 and recombination rates recommended for humans in MSMC's manual, respectively

278 $1.25 \times 10^{-8}$ and $1 \times 10^{-8}$ (https://github.com/stschiff/msmc/blob/master/guide.md).

279 When analyzing simulated sequence data, we simulate sequences of 100 Mb: two se-

280 quences for eSMC and MSMC2, four sequences for MSMC and twenty sequences for

281 SMC++.

### 282 2.2.2 Calculation of the mean square error (MSE)

283 To measure the accuracy of inferences we calculate the Mean Square Error (MSE).

284 We first divide the time window (in log10 scale) of each analysis into ten thousand

285 points. We then calculate the MSE by comparing the actual population size and the

286 one estimated by the method at each of the ten thousand points. We thus have the

287 following formula:

$$MSE = \frac{\sum_{i=1}^{10^4} (y_i - y_i^*)^2}{10^4} \qquad (3)$$

288 Where:

289 • $y_i$ is the population size at the time point $i$.

290 • $y_i^*$ is the estimated population size at the time point $i$.

291 All the command lines to simulate data can be found in S1 of the Appendix.

13

### 2.2.3   Presence of hypothesis violation

**SNP calling:** In practice, SNP calling from next generation sequencing can yield different numbers and frequencies of SNPs depending on the chosen parameters for the different steps of analysis (read trimming, quality check, read mapping, and SNP calling) as well as the quality of the reference genome, data coverage and depth of sequencing, species ploidy [46]. Therefore, based on raw sequence data, the stringency of filters can lead to excluding SNPs (false negatives) or including spurious ones (false positives). When dealing with complex genomes or ancient DNA [56, 7], SNPs can be simultaneously missed and added. We thus simulate four sequences of 100 Mb under a "saw-tooth" scenario and then a certain percentage (5,10 and 25 % ) of SNPs is randomly added to and/or deleted from the simulated sequences. We then analyze the variation and bias in SNP calling on the accuracy of demographic parameter estimations. As an additional analysis we test the effect of ascertainment bias on inferences (a prominent issue in microarray SNP studies) by simulating 100 sequences with msprime where only SNPs above a certain (Minor Allele Frequency) MAF threshold (1%,5% and 10%) are kept, then run SMC methods on a subset of the obtained data.

**Changes in mutation and recombination rates along the sequence:** Because the recombination rate and the mutation rate can change along the sequence [1], and chromosomes are not always fully assembled in the reference genome (which consists of possibly many scaffolds), we simulate short sequences where the recombination and/or mutation rate randomly change between the different scaffolds around an average value of $1.25 \times 10-8$ per generation per base pair (between $2.5 \times 10-9$ and $6.25 \times 10-8$). We simulate 20 scaffolds of size 2 Mb, as this seems representative of

14

316 the best available assembly for non-model organisms [28, 61]. We then analyze the

317 simulated sequences to study the effect of assuming scaffolds share the same muta-

318 tion and recombination rates. In addition, we simulate sequences of 40 Mb (assuming

319 genomes are fully assembled) where the recombination rate along the sequence ran-

320 domly changes every 2 Mbp (up to five-fold) around an average value of $1.25 \times 10{-}8$

321 (the mutation rate being fixed at $1.25 \times 10{-}8$ per generation per bp) to study the

322 effect of the assumption of a constant recombination rate along the sequence.

323 **Transposable elements (TEs):** Genomes can contain transposable ele-

324 ments whose dynamics violate the classic infinite site mutational model for SNPs,

325 and thus potentially affect the estimation of different parameters. Although methods

326 have been developed to detect [39] and simulate them [24], understanding how their

327 presence/absence influences demographic inferences remains unclear. TEs are usually

328 masked when detected in the reference genome and thus not taken into account in the

329 mapped individuals due to the redundancy of read mapping for TEs. Due to their

330 repetitive nature, it can be difficult to correctly detect and assemble them if using

331 short reads, as well as to assess the presence/absence polymorphism of individuals of

332 a population [11]. In addition, the quality and completeness of the reference genome

333 (*e.g.* using the reference genome of a sister species as the reference genome) can

334 strongly affect the accuracy of detecting, assembling and masking TEs [47]. To best

335 capture and mimic the effect of TEs unaccounted for in the data, we altered four sim-

336 ulated sequences of length 20 Mb in four different ways. The first way to simulate the

337 effect of unmapped and unaccounted TEs is to assume they exhibit presence/absence

338 polymorphism, hence creating gaps in the sequence. For each individual, we remove

339 small pieces of sequence of different length (1kb, 10 kb or 100kb), so that up to a
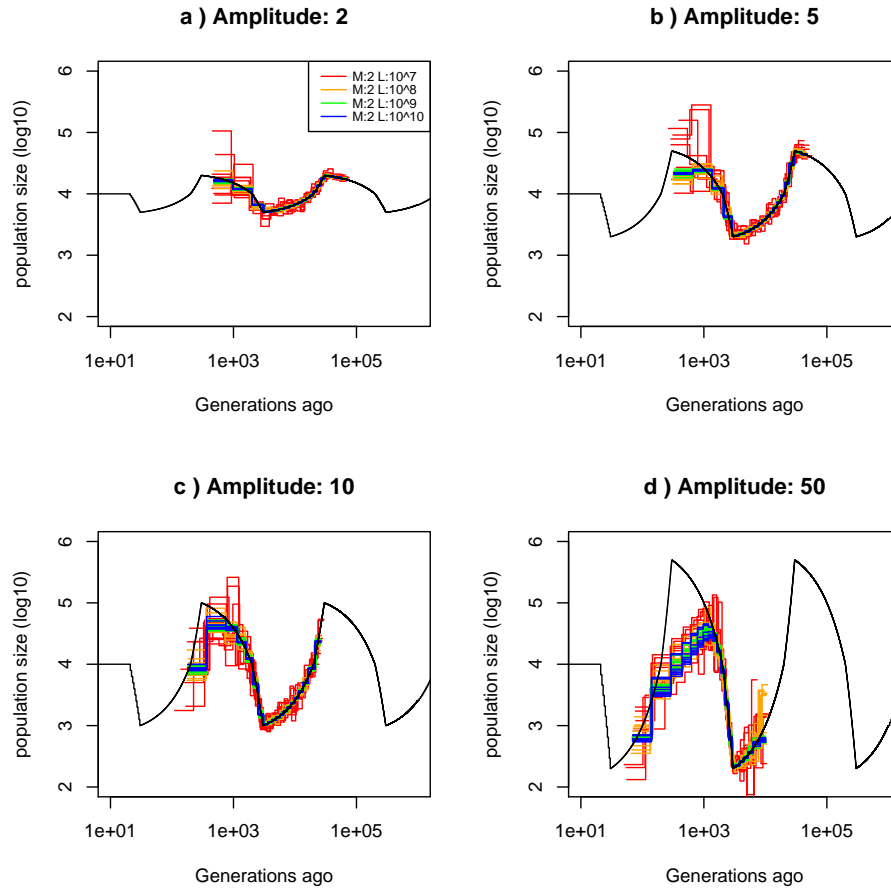
15

340 certain percentage (5,10,25,50%) of the original simulated sequence is removed, so

341 as to shorten and fragment the whole sequence to be analyzed. The second way, is

342 to consider unmasked TEs, done by randomly selecting small pieces of the original

343 simulated sequence (1kb, 10 kb or 100kb), making up to a certain percentage of it

344 (5,10,25,50%), and removing all the SNPs found in those regions (*i.e.* removing muta-

345 tions from TEs). The removed SNPs are hence structured in many small regions along

346 the genome. Thirdly, we test the consequences of simultaneously having both removed

347 and unmasked TEs in the data set. Lastly, to measure the importance of detecting

348 and masking TEs, we assume all TEs to be present and masked when building the

349 multihetsep file (*i.e.* considering TEs as missing data).

## 3  Results

### 3.1  Best-case convergence

352 Results of the best-case convergence of eSMC under the saw-tooth demographic his-

353 tory are displayed in Figure 1. Increasing the sequence length increases accuracy and

354 reduces variability, leading to better convergence and reducing the mean square error

355 (see Figures 1a-c and Supplementary Table 1). However, when the amplitude of popu-

356 lation size variation is too great (here for 50 fold), the demographic history cannot be

357 retrieved, even when using very large data sets (see Figure 1d). Similar results are ob-

358 tained for the three other demographic scenarios (bottleneck, expansion and decrease,

359 respectively displayed in Supplementary Figures 1, 2 and 3). The bottleneck scenario

360 seems especially difficult to infer, requiring large amounts of data, and the stronger

361 the bottleneck, the harder it is to detect it, even with sequence lengths equivalent to

362 $10^{10}$bp. In Supplementary Figure 4, we show that even when changing the number of

16

363    hidden states (*i.e.* number of inferred parameters), some scenarios with very strong
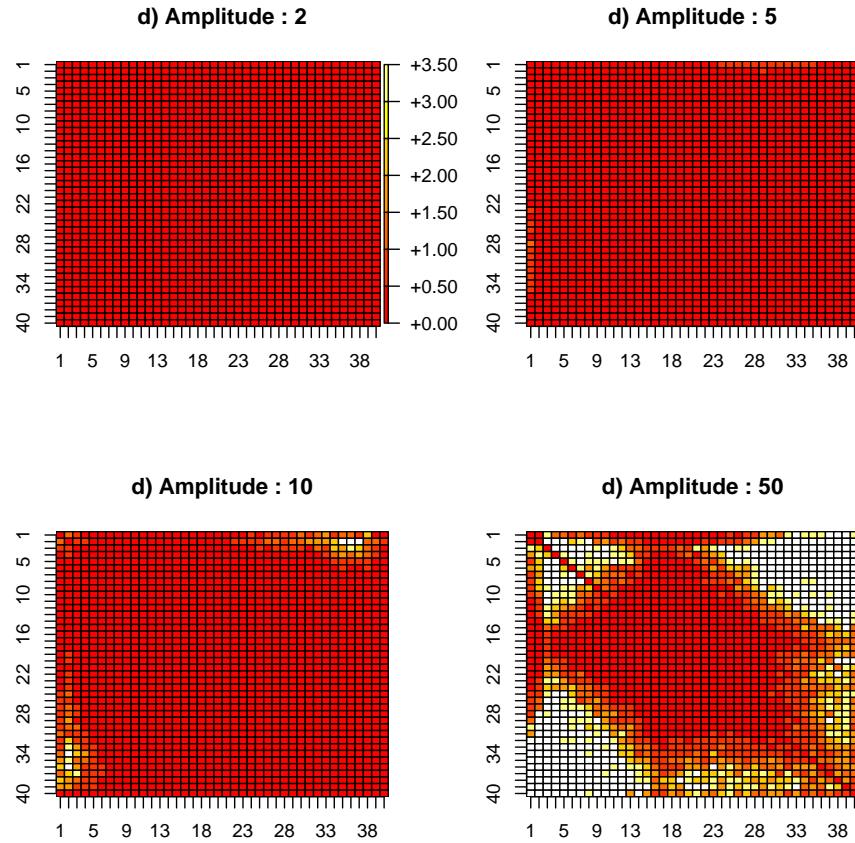
364    variation of population size remain badly inferred.

**Fig. 1.** **Best-case convergence of eSMC** Estimated demographic history using simulated genealogy over sequences of 10,100,1000,10000 Mb (respectively in red,orange, green and blue) under a saw-tooth scenario (original scenario in black) with 10 replicates for different amplitudes of size change: a) 2-fold, b) 5-fold, c) 10-fold, and d) 50-fold. The recombination rate is set to $1 \times 10^{-8}$ per generation per bp and the mutation rate to $1.25 \times 10^{-8}$ per generation per bp.

**365**     In Supplementary Figures 5, 6, 7 and 8, we show the best-case convergence of

18

366 MSMC with four genome sequences and generally find that these analyses present a

367 higher variance than eSMC. However, MSMC shows better fits in recent times and is

368 better able to retrieve population size variation than eSMC (see Supplementary Figure

369 5d). Scenarios with strong variation of population size (*i.e.* with large amplitudes) still

370 pose a problem (see Supplementary Figure 9), and no matter the number of estimated

371 parameters, such scenarios cannot be correctly inferred using MSMC.

372 To better understand these results, we examine the coefficient of variation cal-

373 culated from the replicates at each entry of the transition matrix. We can see that

374 increasing the sequence length reduces the coefficient of variation (the ratio of the

375 standard deviation to the mean, hence indicating convergence when equal to 0, see

376 Supplementary Figure 10). Yet increasing the amplitude of population size variation

377 decreases the number of some hidden state transitions leading to unobserved transi-

378 tions. Unobserved transitions result from the reduced probability of coalescence events

379 in specific time intervals (*i.e.* hidden states). In these cases matrices display higher co-

380 efficients of variation and can be partially empty (Figure 2). This explains the increase

381 of variability of the inferred scenarios, as well as the incapacity of SMC methods to

382 correctly infer the demographic history with strong population size variation in specific

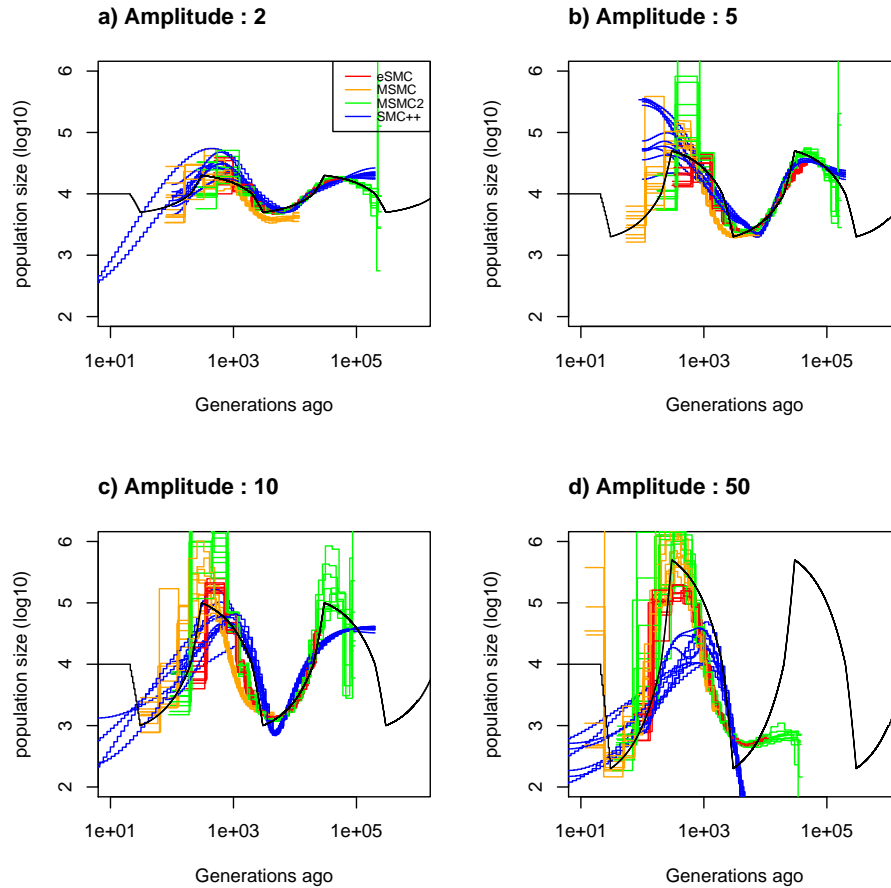383 time intervals independently of the amount of data available.

19

**Fig. 2. Estimated transition matrix in sharp saw-tooth scenario** Estimated coefficient of variation of the transition matrix using simulated genealogy over sequences of 10000 Mb under a saw-tooth scenario of amplitude 2, 5,10 and 50 (respectively in a, b, c and d) each with 10 replicates. Recombination and mutation rates are as in Figure 1. White squares indicate absence of observed transitions (*i.e. no data*).

20

## 3.2 Simulated sequence results

### 3.2.1 Scenario effect

In the previous section, we explored the theoretical performance limitations of eSMC and MSMC using trees in Newick format as input. In this section, we evaluate how these methods perform when inputting simulated sequence data using the same recombination and mutation rates. We first perform two benchmark analyses, the constant population size scenario (Supplementary Figure 11) and the sawtooth demographic scenario from [50] (Supplementary Figure 12). eSMC and MSMC2 retrieve the constant population size scenario although MSMC fails to retrieve it in the far past and SMC++ in recent time (Supplementary Figure 11). All methods can retrieve the sawtooth demographic scenario although SMC++ displays high variance in recent times (Supplementary Figure 12). Secondly, we investigate the effect of amplitude of population size variation as in Figure 1. Results for the saw-tooth scenario are displayed in Figure 3, where the different models display a good fit, but are not as good as expected from the best-case convergence given the same amount of data (Figure 1 (orange line) and Supplementary Table 1 vs Figure 3 (red line) and Supplementary Table 2). As predicted by Figures 1 and 2, the case with the greatest amplitude of population size variation (Figure 1d) is the least well fitted (see Supplementary Table 2 for the MSE). All estimations display low variance and a relatively good fit in the bottleneck and expansion scenarios for small population size variation (see Supplementary Figures 13a and 14a ). However, the strengths of expansions and bottlenecks are not fully retrieved in scenarios with population size variation equals to or is higher than tenfold the current population size (Supplementary Figures 13c-d,and 14c-d). To study the origin of differences between simulation results and theoretical results, we measure

21

**408** the difference between the transition matrix estimated by eSMC and the one built

**409** from the actual genealogy. Results show that hidden states are harder to correctly

**410** infer in scenarios with strong population size variation, explaining the high variance

**411** (see Supplementary Figure 15). We demonstrate there that for the same amount of

**412** data, the simulation, and thus by extension the real data, shows additional stochastic

**413** behaviour than the best-case convergence (Figure 1).

**Fig. 3. Estimated demography using simulated sequences as input.** Estimated demographic history (black) under a saw-tooth scenario with 10 replicates using simulated sequences for different amplitudes of population size change: a) 2, b) 5, c) 10 and d) 50. Two sequences of 100 Mb for eSMC and MSMC2 (respectively in red and green), four sequences of 100 Mb for MSMC (orange) and 20 sequences of 100 Mb for SMC++ (blue) were simulated. Recombination and mutation rates are respectively set to $1 \times 10^{-8}$ and $1.25 \times 10^{-8}$.

23

**414**       Increasing the time window in eSMC results in an increased variance of the

**415** inferences (Supplementary Figure 16). In addition, shifting the window towards more

**416** recent time leads to poor demographic estimations, but shifting it further in the past

**417** does not seem to bias it (there are however consequences on estimations of the recom-

**418** bination rates, see Table 1 for more details). Concerning the optimization function, we

**419** find that the complete Baum-Welch algorithm gives similar results to the incomplete
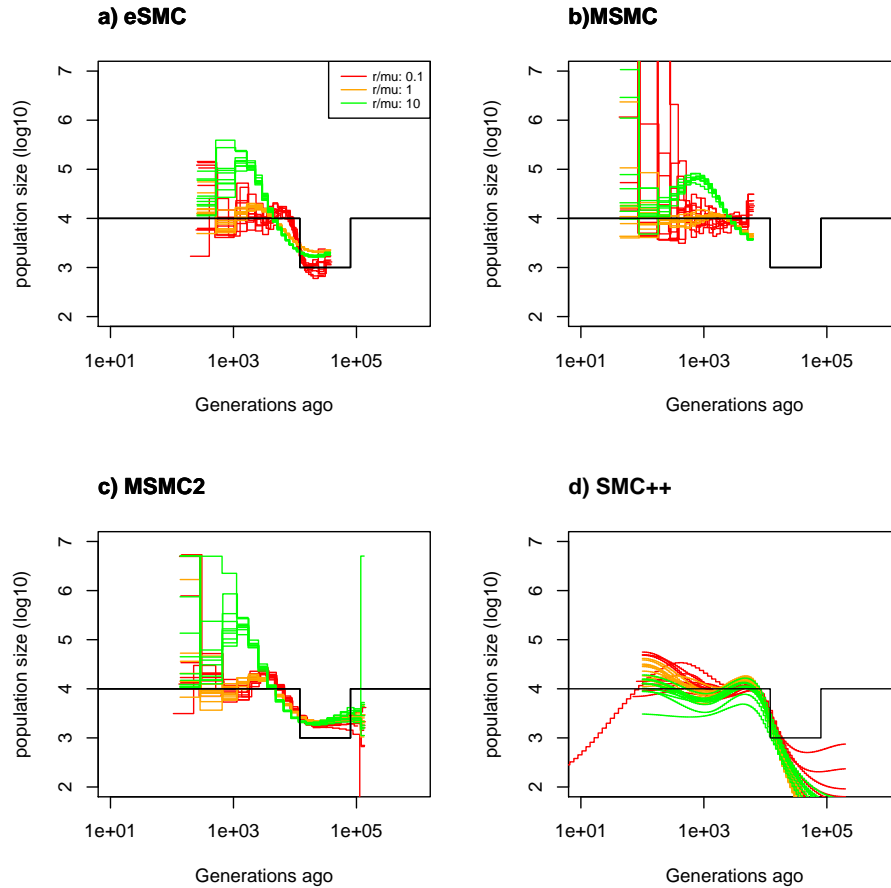
**420** one (Table 1).

| Optimization function | Scenario | real $\frac{\rho}{\theta}$ | normal window $\frac{\rho}{\theta}^*$ | Big Window $\frac{\rho}{\theta}^*$ | Old window $\frac{\rho}{\theta}^*$ | Recent window $\frac{\rho}{\theta}^*$ |
|---|---|---|---|---|---|---|
| Incomplete Baum-Welch | Saw-tooth | 0.8 | 0.79 (0.036) | 0.72 (0.039) | 0.72 (0.042) | 0.94 (0.005) |
| Complete Baum-Welch | Saw-tooth | 0.8 | .79 (0.044) | 0.72 (0.039) | 0.72 (0.042) | 1.56 (0.087) |
| Incomplete Baum-Welch | Constant | 0.8 | 0.86 (0.019) | 0.85 (0.020) | 0.84 (0.019) | 0.98 (0.002) |
| Complete Baum-Welch | Constant | 0.8 | 0.86 (0.019) | 0.85 (0.020) | 0.84 (0.019) | 1.06 (0.02) |

Table 1: Average estimated values for the recombination over mutation ratio $\frac{\rho}{\theta}$ by eSMC over ten repetitions for different sizes of the time window. The coefficient of variation is indicated in brackets. Four sequences of 50 Mb were simulated with a recombination rate set to $1 \times 10^{-8}$ per generation per bp and a mutation rate to $1.25 \times 10^{-8}$ per generation per bp.

**421**       Adding a regularization penalty to eSMC can drastically impact inferences

**422** (Supplementary Figure 17) and reduces performance quality. When regularization is

**423** added, eSMC fails to correctly capture the amplitude of population size variation and

**424** with extreme regularization penalty, eSMC infers a constant population size. Yet,

**425** adding regularization in SMC++ can increase performance and avoid spurious vari-

**426** ation of population size (Supplementary Figure 18). However, strong regularization

**427** can lead to the inference of constant population size and thus poor estimations.

24

### 3.2.2 Effect of the ratio of the recombination over the mutation rate

The ratio of the effective recombination over effective mutation rates ($\frac{\rho}{\theta}$) can influence the ability of SMC-based methods to retrieve the coalescence time between two points along the genome [63]. Intuitively, if recombination occurs at a higher rate compared to mutation, then it renders it more difficult to detect any recombination events that may have taken place before the introduction of a new mutation, and thus bias the estimation of the coalescence time [53, 63]. Under the bottleneck scenario, we find that the lower $\frac{\rho}{\theta}$, the better the fit of the inferred demography by eSMC and SMC++ in the past, but also the higher the variance of the inferences (see Figure 4 ). However each method displays the worse fit when $\frac{\rho}{\theta} = 10$ (Supplementary Table 3). SMC++ seems slightly less sensitive to $\frac{\rho}{\theta}$ than other methods. When calculating the difference between the transition matrix estimated by eSMC and the one built from the actual genealogy (ARG), we find that, unsurprisingly, changes in hidden states are harder to detect when $\frac{\rho}{\theta}$ increases, leading to an overestimation of hidden states on the diagonal (see Supplementary Figures 19, 20 and 21).

**Fig. 4. Effect of $\frac{\rho}{\theta}$ on inference of demographic history.** Estimated demographic history under a bottleneck scenario with 10 replicates using simulated sequences. We simulate two sequences of 100 Mb for eSMC and MSMC2 (respectively in a and b), four sequences of 100 Mb for MSMC (c) and twenty sequences of 100 Mb for SMC++ (d). The mutation rate is set to $1.25 \times 10^{-8}$ per generation per bp and the recombination rates are $1.25 \times 10^{-9}, 1.25 \times 10^{-8}$ and $1.25 \times 10^{-7}$ per generation per bp, giving $\frac{\rho}{\theta} = 0.1$, 1 and 2 and the inferred demographies are in red, orange and green respectively. The demographic history is simulated under a bottleneck scenario of amplitude 10 and is represented in black.

26

**443** It is, in some instances, possible to compensate for a $\frac{\rho}{\theta}$ ratio that is not ideal

**444** by increasing the number of iterations. Indeed, for eSMC, the demographic history

**445** is better inferred (Supplementary Figure 22), although the correct recombination rate

**446** cannot be retrieved (Table 2). MSMC is able to better infer the correct recombina-

**447** tion rate than other methods despite $\frac{\rho}{\theta} > 1$, but poorly estimates the demographic

**448** history. The past demographic inferences obtained using MSMC2 and SMC++ are

**449** not improved when increasing the number of iterations (see Supplementary Figure 22

**450** and Table 2).

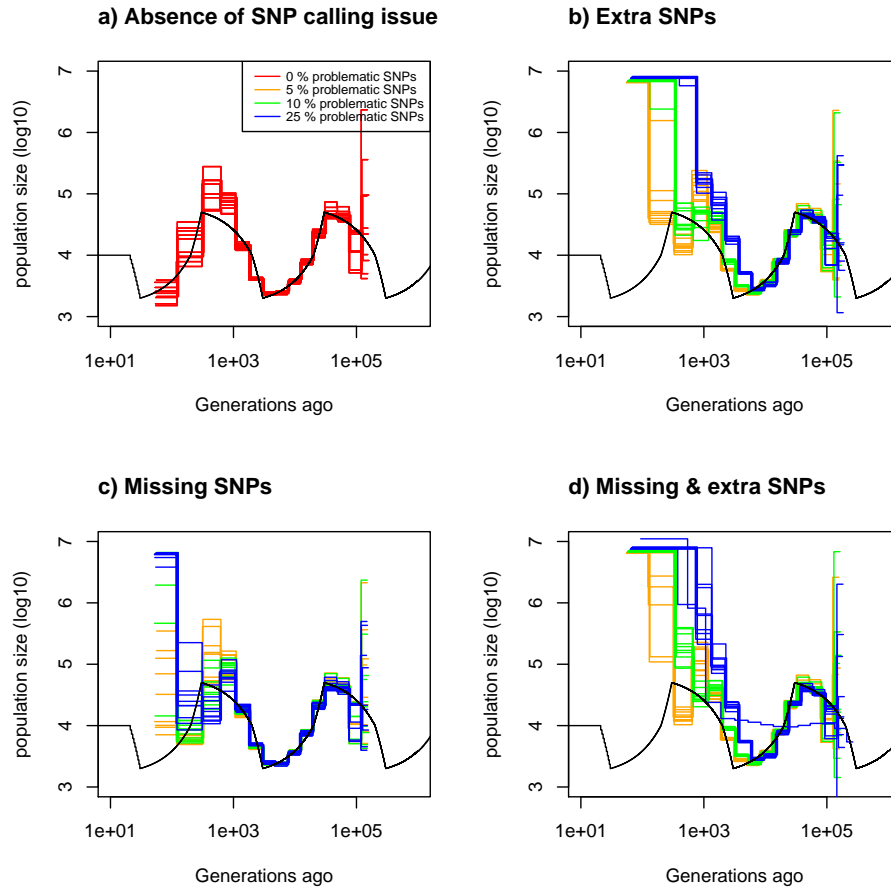| method | real $\frac{\rho}{\theta}$ | set 1 , $\frac{\rho}{\theta}^*$ | set 2 , $\frac{\rho}{\theta}^*$ | set 3 , $\frac{\rho}{\theta}^*$ | set 4 , $\frac{\rho}{\theta}^*$ | set 5 , $\frac{\rho}{\theta}^*$ |
|--------|------|------|------|------|------|------|
| eSMC | 10 | 1.35 (0.026) | 1.76 (0.047) | 1.29 (0.027) | 1.74 (0.048) | 1.80 (0.041) |
| MSMC | 10 | 2.70 (0.011) | 6.58 (0.031) | 2.68 (0.011) | 6.57 (0.032) | 6.62 (0.030) |
| MSMC2 | 10 | 1.27 (0.055) | 1.65 (0.13) | 1.26 (0.060) | 1.75 (0.060) | 1.60 (0.29) |
| SMC++ | 10 | 0.56 (0.38) | 0.48 (0.38) | 1.32 (0.15) | 0.21 (0.62) | 0.98 (0.24) |

Table 2: Average estimated values for the recombination over mutation ratio $\frac{\rho}{\theta}$ over ten repetitions. The coefficient of variation is indicated in brackets. For eSMC, MSMC and MSMC2 we have: set 1: 20 hidden states; set 2: 200 iterations; set3: 60 hidden states; set 4: 60 hidden states and 200 iterations and set 5: 20 hidden states and 200 iterations. For SMC++: set 1: 16 knots; set 2: 200 iterations; set 3: 4 knots in green; set 4: regularization penalty set to 3 and set 5: regularization-penalty set to 12.

**451** ## 3.3 Simulation results under hypothesis violation

**452** ### 3.3.1 Imperfect SNP calling

**453** We analyze simulated sequences that have been modified by removing and/or adding

**454** SNPs using the different SMC methods. We find that, when using MSMC2, eSMC

455   and MSMC, having more than 10% of spurious SNPs (*e.g.* no quality filtering) can

456   lead to a strong over-estimation of population size in recent time but that missing

457   SNPs have no effects on inferences in the far past and only mild effects on inferences

458   of recent time (see Figure 5 for MSMC2, Supplementary Figures 23 and 24 for eSMC

459   and MSMC respectively). The mean square error is displayed in Supplementary Table
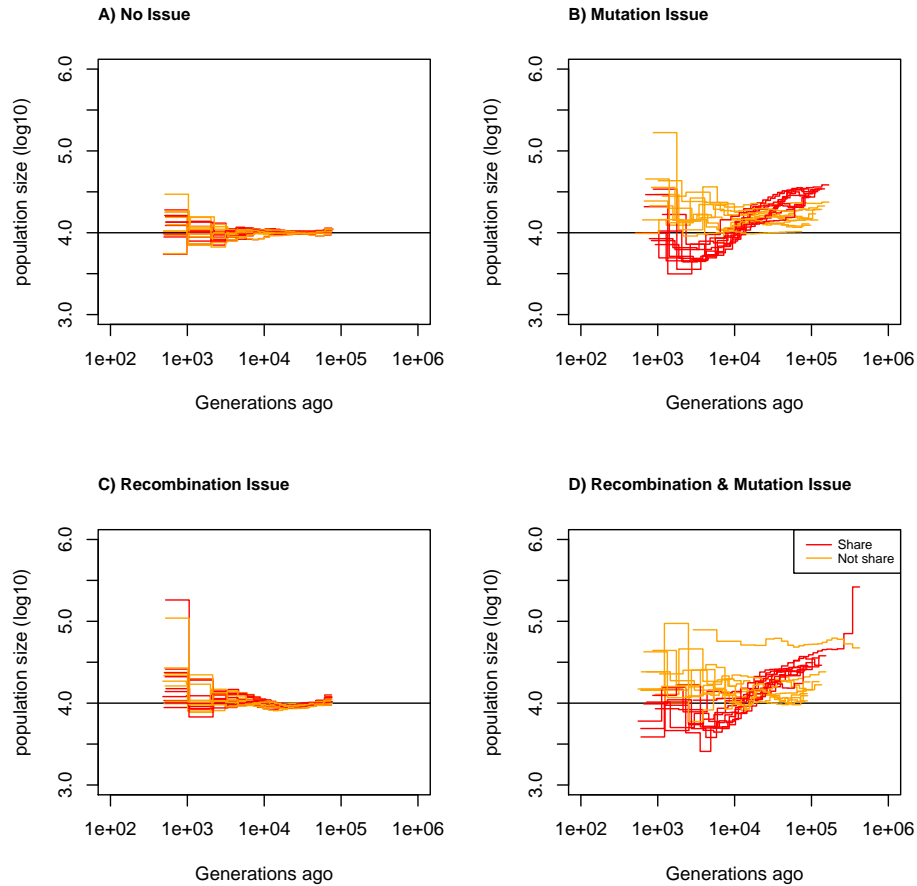
460   4.

**Fig. 5. Consequences of SNP calling errors.** Estimated demographic history using MSMC2 under a saw-tooth scenario with 10 replicates using four simulated sequences of 100 Mb. Recombination and mutation rates are as in Figure 1 and the simulated demographic history is represented in black. a) Demographic history simulated with ibsence of SNP calling issue (red). b) Demographic history simulated with 5% (orange),10% (green) and 25% (blue) missing SNPs. c) Demographic history simulated with 5% (orange), 10% (green) and 25% (blue) additional SNPs. d) Demographic history simulated with 5% (orange),10% (green) and 25% (blue) of additional and missing SNPs .

29

461     As complementary analyses we analyze simulated sequences with a Minor Allele

462 Frequency (MAF) threshold. Results are shown in Supplementary Figure 25. The

463 more SNPs are removed, the poorer the estimations in recent time (Supplementary

464 Figure 25), demonstrating the impact of severe ascertainment bias.

465 **3.3.2   Specific scaffold parameters**

466 We simulate sequence data where scaffolds have either been simulated with the same

467 recombination and mutation rates or with different recombination and mutation rates.

468 Data sets are then analyzed assuming scaffolds share or do not share the same re-

469 combination and mutation rates. We can see in Figure 6 (and Supplementary Table

470 5) that when scaffolds all share the same parameter values, estimated demography is

471 accurate both when the analysis assumed shared or differing mutation and recombi-

472 nation rates. However, when scaffolds are simulated with different parameter values,

473 analyzing them under the assumption that they have the same mutation and recombi-

474 nation rates leads to poor estimations. Assuming scaffolds do not share recombination

475 and mutation rates does improve the results somewhat, but the estimations remain

476 less accurate than when scaffolds all share with same parameter values. If only the

477 recombination rate changes from one scaffold to another, the demographic history is

478 only slightly biased, whereas, if the mutation rate changes from one scaffold to the

479 other, demographic history is poorly estimated.

30

**Fig. 6. Estimating demographic history using scaffolds sharing or differing in mutation and recombination rates** Estimated demographic history using eSMC under a saw-tooth scenario with 10 replicates using twenty simulated scaffolds of two sequences of 2 Mb assuming scaffolds share (red) or do not share recombination and mutation rates (orange). The simulated demographic history is represented in black. a) Scaffolds share the same parameters, recombination and mutation rates are set at $1.25 \times 10^{-8}$, b) Each scaffold is randomly assigned a recombination rate between $2.5 \times 10^{-9}$ and $6.25 \times 10^{-8}$ and the mutation rate is $1.25 \times 10^{-8}$, c) Each scaffold is randomly assigned a mutation rate between $2.5 \times 10^{-9}$ and $6.25 \times 10^{-8}$ and the recombination rate is $1.25 \times 10^{-8}$ and d) Each scaffold is assigned a random mutation and an independently random recombination rate, both being between $2.5 \times 10^{-9}$ and $6.25 \times 10^{-8}$.
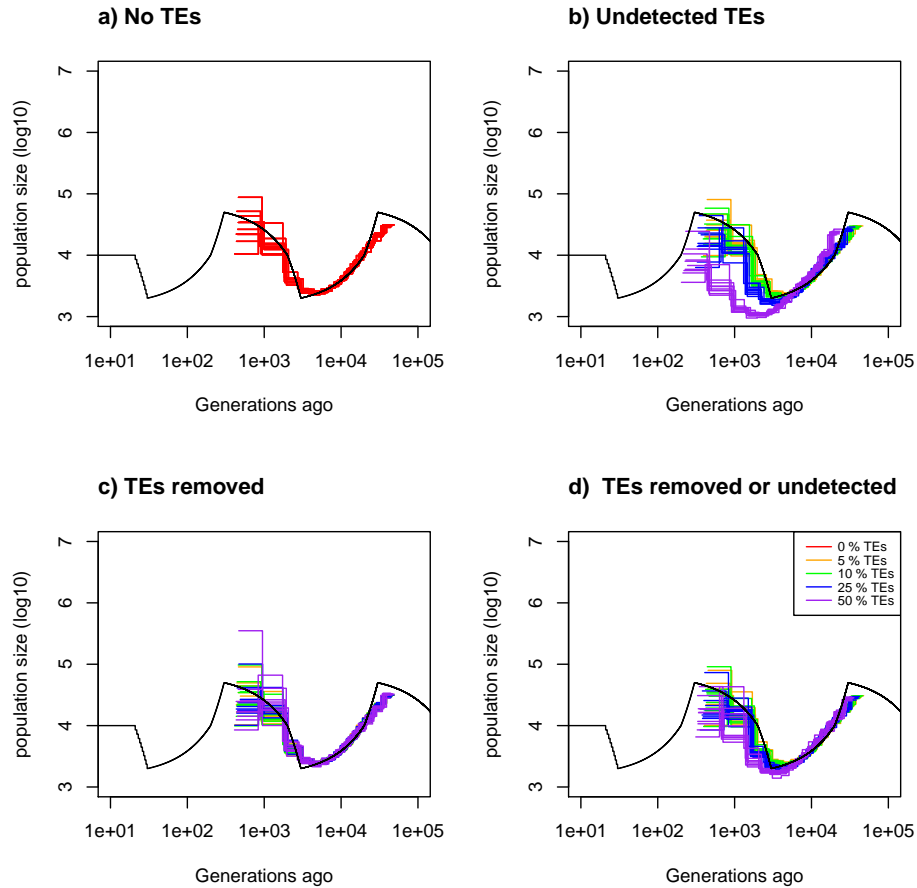
34

480    Even if chromosomes are fully assembled, assuming we here have one scaffold

481    of 40 Mb (chromosome fully assembled), there may be variations of the recombination

482    rate along the sequence, however this seems of little consequence when applying eSMC.

483    As can be seen in Supplementary Figure 26, the demographic scenario is well inferred,

484    despite an increase in variance and a smooth "wave" shaped demographic history

485    when sequences simulated with varying recombination rates are compared to those

486    with a fixed recombination rate throughout the genome. Overall we see that when

487    recombination rate is heterogeneous along the genome by a factor 5, it is not untypical

488    to falsely estimate a two-fold variation of Ne even though the true Ne is constant in

489    time.


### 3.3.3   How transposable elements bias inference

491    Transposable elements (TEs) are present in most species, and are (if detected) taken

492    into account as missing data by SMC methods [50]). Depending on how TEs affect

493    the data set, we find that methods are more or less sensitive to TEs. If TEs are

494    unmapped/removed from the data set, there does not appear to be any bias in the

495    estimated demographic history when using eSMC (see Figure 7 and Supplementary

496    Table 6), but there is an overestimation of $\frac{\rho}{\theta}$ (see Table 3). We find that, the higher the

497    proportion of sequences removed, the more $\frac{\rho}{\theta}$ is over-estimated. For a fixed amount

498    of missing/removed data, the smaller the sequences that are removed, the more $\frac{\rho}{\theta}$ is

499    over-estimated (Table 3). If TEs are present but unmasked in the data set (and thus

500    not accounted for missing data by the model [50] ), we find that this is equivalent to

501    a faulty calling of SNPs, in which SNPs are missing, hence resulting in demographic

502    history estimations by eSMC similar to those observed in Figure 5a. However, if

503    the size of unmasked TEs increases, different results are obtained (see Supplementary

32

504  Figures 27 and 28). Indeed, in recent times there is a strong underestimation of

505  population size and the model fails to capture the correct demographic history. The

506  longer the TEs are, the stronger the effect on the estimated demographic history.

507  Similar results are obtained with MSMC (Supplementary Figures 29, 30 and 31) and

508  MSMC2 (Supplementary Figures 32, 33 and 34). However, when TEs are detected

509  and correctly masked, there is no effect on demographic inferences (Supplementary

510  Figures 35 and 36).

**Fig. 7. Consequences of masking or removing transposable elements (TEs) from data sets.** Estimated demographic history by eSMC under a saw-tooth scenario with 10 replicates using four simulated sequences of 20 Mb. The recombination and mutation rates are as in Figure 1 and the simulated demographic history is represented in black. Here the TEs are of length 1kbp. a) Demographic history simulated with no TEs. b) Demographic history simulated where TEts are removed. c) Demographic history simulated where TEs are masked. d) Demographic history simulated where half of the TEs are removed and SNPs on the other half are removed. Proportion of the genome made up by TEs is set to 0% (red), 5% (orange), 10% (green), 25 % (blue) and 50 % (purple).

34

| TE length | method | real $\frac{\rho}{\theta}$ | $\frac{\rho}{\theta}^*$ and 5% TEs | $\frac{\rho}{\theta}^*$ and 10% TEs | $\frac{\rho}{\theta}^*$ and 25% TEs | $\frac{\rho}{\theta}^*$ and 50% TEs |
|---|---|---|---|---|---|---|
| | eSMC | 1 | 0.95 (0.021) | 0.99 (0.022) | 1.16 (0.10) | 1.77 (0.36) |
| 1 kb | MSMC | 1 | 1.31 (0.098) | 1.35 (0.11) | 1.50 (0.088) | 1.91 (0.11) |
| | MSMC2 | 1 | 0.87 (0.047) | 0.88 (0.049) | 1.0 (0.036) | 1.35 (0.035) |
| | eSMC | 1 | 0.96 (0.053) | 0.98 (0.066) | 1.10 (0.18) | 1.36 (0.41) |
| 10 kb | MSMC | 1 | 1.38 (0.074) | 1.41 (0.0.090) | 1.54 (0.11) | 1.68 (0.13) |
| | MSMC2 | 1 | 0.87 (0.064) | 0.89 (0.067) | .99 (0.15) | 1.13 (0.30) |
| | eSMC | 1 | 0.95 (0.047) | 0.95 (0.051) | 0.98 (0.070) | 1.0 (0.12) |
| 100 kb | MSMC | 1 | 1.36 (0.048) | 1.36 (0.062) | 1.40 (0.093) | 1.49 (0.12) |
| | MSMC2 | 1 | 0.87 (0.056) | 0.88 (0.050) | 0.91 (0.079) | 0.91 (0.073) |

Table 3: Average estimated values for the recombination over mutation ratio $\frac{\rho}{\theta}$ over ten repetitions. The coefficient of variation is indicated in brackets. TEs are of length 1kb, 10kb or 100 kb and are completely removed and the proportion of the genome made up by TEs is 5%,10% ,25% and 50%.

# 4    Discussion

Throughout this work we have outlined the limits of PSMC' and MSMC methodologies, which had, until now, not been clearly defined. We find that, in most cases, if enough genealogies (*i.e.* data) are inputted then the demographic history is accurately estimated, tending to results obtained previously [13, 8], however, we find that the amount of data required for an accurate fit depends on the underlying demographic scenario. The differences with previous works stems from estimations being made using the actual series of coalescence times [13, 8], whereas we use the series of hidden states built from the discretization of time summarized in a simple matrix. We also find that some scenarios are better retrieved when using either MSMC or methods based on PSMC', indicating that there are complementary convergence properties

35

522 between these methodologies.

523    We develop a method to indicate if the amount of data is enough to retrieve
524 a specific scenario, notably by calculating the coefficient of variation of the transition
525 matrix using either real or simulated data, and therefore offer guidelines to build
526 appropriate data sets (see also Supplementary Figure 8). Our approach can also be
527 used to infer demographic history given an ARG (using trees in Newick format or
528 sequences of coalescence events), independently of how the ARG has been estimated.
529 Our results suggest that whole genome polymorphism data can be summarized in
530 a transition matrix based on the SMC theory to estimate demographic history of
531 panmitic populations. As new methods can infer genealogies better and faster [58, 22,
532 36, 42], the estimated transition matrix could become a powerful summary statistic
533 in the future. HMM can be a computational burden depending on the model and
534 model parameters, and estimating genealogy through more efficient methods would
535 still allow the use of SMC theory for parameter estimation or hypothesis testing (as
536 in [67, 13, 19]). In addition, using the work of [66], one could (to some extent [23])
537 extend our approach to account for population structure and migration.

538    We have also demonstrated that the power of PSMC', MSMC, and other SMC-
539 based methods, rely on their ability to correctly infer the genealogies along the se-
540 quence (*i.e.* the Ancestral Recombination Graph or ARG). The accuracy of ARG
541 inference by SMC methods, however, depends on the ratio of the recombination over
542 the mutation rate ($\frac{\rho}{\theta}$). As this rate increases, estimations lose accuracy. Specifically,
543 increasing $\frac{\rho}{\theta}$ leads to an over-estimation of transitions on the diagonal, which explains
544 the underestimation of the recombination rate and inaccurate demographic history es-
545 timations, as shown in [63, 53]. As a way around this issue, in some cases it is possible

36

546   to obtain better results by increasing the number of iterations. MSMC's demographic

547   inference is more sensitive to $\frac{\rho}{\theta}$ but the quality of the estimation of the ratio itself is

548   less affected. This once again shows the complementarity of PSMC' and MSMC. If

549   the variable of interest is $\frac{\rho}{\theta}$, then MSMC should be used, but if the demographic his-

550   tory is of greater importance, PSMC'-based methods should be used. The amplitude

551   of population size variation also influences the estimation of hidden states along the

552   sequences, with high amplitudes leading to a poor estimation of the transition ma-

553   trix, distorting the inferred demography. We find that increasing the size of the time

554   window increases the variance of the estimations, despite using the same number of

555   parameters, as this results in a small under-estimation of $\frac{\rho}{\theta}$. In addition the complete

556   and incomplete Baum-Welch algorithms lead to identical results, demonstrating that

557   all the information required for the inference is in the estimated transition matrix.

558   Finally, we explored how imperfect data sets (due to errors in SNP calling,

559   the presence of transposable elements and existing variation in recombination and

560   mutation rates) could affect the inferences obtained using SMC-based methods. We

561   show that a data set with more than 10% of spurious SNPs will lead to poor estimations

562   of the demographic history, whereas randomly removed SNPs (*i.e.* missing SNPs) have

563   a lesser effect on inferences. It is thus better to be stringent during SNP calling, as

564   SNPs is worse than missing SNPs. Note, however, that this consideration is valid for

565   demographic inference under a neutral model of evolution, while biases in SNP calling

566   also affect the inference of selection (especially for conserved genes under purifying

567   selection). However, if missing SNPs are structured along the sequence (as would be

568   the case with unmasked TEs), there is a strong effect on inference. If TEs are correctly

569   detected and masked, there is no effect on demographic inferences. It is therefore

37

570 recommended that checks should be run to detect regions with abnormal distributions

571 of SNPs along the genome. Surprisingly, simulation results suggest that removing

572 random pieces of sequences have no impact on the estimated demographic history.

573 Taking this into account, when seeking to infer demographic history, it seems better

574 to remove sections of sequences than to introduce sequences with SNP call errors or

575 abnormal SNP distributions. However, removing sequences leads to an over-estimation

576 of $\frac{\rho}{\theta}$, which seems to depend on the number and size of the removed sections. The

577 removal of a few, albeit long sequences, will have almost no impact, whereas removing

578 many short sections of the sequences will lead to a large overestimation of $\frac{\rho}{\theta}$. This

579 consequence could provide an explanation for the frequent overestimation of $\frac{\rho}{\theta}$ when

580 compared to empirical measures of the ratio of recombination and mutation rates $\frac{r}{\mu}$.

581 This implies, that in some cases, despite an inferred $\frac{\rho}{\theta} > 1$, the inferred demographic

582 history can surprisingly be trusted. Note also that as discussed in [53], the discrepancy

583 between $\frac{\rho}{\theta}$ and $\frac{r}{\mu}$ can be due to life history traits such as selfing or dormancy.

584 Simulation results suggest that any variation of the recombination rate along

585 the sequence does not strongly bias demographic inference but slightly increases the

586 variance of the results and leads to small waves in the demographic history (as a

587 consequence of erroneously estimated hidden state transition events because of the

588 non constant recombination rate along the sequence), as expected from previous works

589 [26]. However, unlike Li and Durbin's results [26], if scaffolds do not share similar

590 rates of mutation and recombination, but are analyzed together assuming that they

591 do, estimations will be very poor. This could be due to the variation of mutation

592 rate being within a scaffold in their study and the discrepancy between out and their

593 results could suggest analyses based on longer scaffolds to be more robust. However,

38

594 this problem can be avoided if each scaffold is assumed to have its own parameter

595 values, although this would increase computation time, it could provide useful insight

596 in unveiling any variation in molecular forces along the genome, albeit in a coarser

597 way than in [1]. As we found that non-accounted variation of the recombination rate

598 along the sequence can lead to a spurious two-fold variation of population size, we

599 here provide guidelines to test if small detected variations of population size are to be

600 trusted. Since the consequecnes of a varying recombination rate might depend on the

601 topology of the recombination map, one first needs estimate the recombination map

602 (*e.g.* using iSMC [1]). If problematic regions are found they can be removed with

603 almost no negative impact on the estimated demography (Figure 7). Otherwise,the

604 recombination map can be used to simulate sequences *e.g.* using scrm [60]), which can

605 be compared to results obtained for a constant recombination rate. Analyses can be

606 run on both data sets to quantify the effect of the recombination map.

## 607 4.1 Guidelines when applying SMC-based methods

608 Our aim through this work is to provide guidelines to optimize the use of SMC-based

609 methods for inference. First, if the data set is not yet built, but there is some intuition

610 concerning the demographic history and knowledge of some genomic properties of a

611 species (*e.g.* recombination and mutation rates), we recommend simulating a data

612 set corresponding to the potential scenarios. From these simulations, the transition

613 matrix for PSMC' or MSMC-based methods can be built using the R package eSMC2.

614 The results obtained can guide users when it comes to the amount and quality of data

615 needed (sequence size and copy number) for a good inference. Beyond being used

616 to guide the building of data sets, it is possible to assess trustworthiness of results

617 obtained using SMC-based methods on existing data sets. If the estimated transition

39

618 matrix is empty in some places (*i.e.* no observed transition event between two specific

619 hidden states; white squares in Figure 2), it could suggest a lack of data and/or strong

620 variation of the population size somewhere in time. In order to test the accuracy of the

621 inferred demography, the estimated demographic history can be retrieved and used to

622 simulate a data set with more sequences and/or simulate a demographic history with

623 a higher amplitude than the estimated one. The SMC method can then be run on

624 the simulated data in order to check whether using more data results in a matching

625 scenario or if a higher amplitude of population size can indeed be inferred, in which

626 cases the initial results are most probably trustworthy.

627 As mentioned above, it is better to sequence fewer individuals, but to have

628 data of better quality. It is also important to note, that more data is not necessarily

629 always better, especially if there is a risk of spurious SNPs (see Figure 5). In some

630 cases, methods are limited by their own theoretical framework, hence no given data set

631 will ever allow a correct demographic inference. In such cases, other methods based on

632 a different theoretical frameworks (*e.g.* SFS and ABC) might perform better [3, 51].

## 4.2 Concluding remarks

634 Here we present a simple method to help assess how accurate inferences obtained us-

635 ing PSMC' and MSMC would be when applied to data sets with suspected flaws or

636 limitations. We also provide new interpretations of results obtained when hypotheses

637 are known to be violated, and thus offer an explanation as to why results sometimes

638 deviate from expectations (*e.g.* when the estimated ratio of recombination over mu-

639 tation is larger than the one measured experimentally). We propose guidelines for

640 building/evaluating data sets when using SMC-based models, as well as a method

40

which can be used to estimate the demographic history and recombination rate given a genealogy (in the same spirit as Popsicle [13]). The estimated transition matrix is introduced as a summary statistic, which can be used to recover demographic history (more precisely the Inverse Instantaneous Coalescence Rate interpretation of population size variation, when assuming a panmictic population [8, 49]). This statistic could, in future, be used in scenarios with migration, without the computational load of Hidden Markov models. When faced with complex demographic histories, or $\frac{\rho}{\theta} > 1$, we show that there are strategies that would allow those wishing to use SMC methodology to make the best use of their data.

# 5   Acknowledgments

# 6   Competing Interest

The authors of this article declare that they have no financial conflict of interest with the content of this article.

# References

[1] Gustavo V. Barroso, Natasa Puzovic, and Julien Y. Dutheil. Inference of recombination maps from a single pair of genomes and its application to ancient samples. *PLOS GENETICS*, 15(11), NOV 2019.

[2]  Champak R. Beeravolu, Michael J. Hickerson, Laurent A. F. Frantz, and Konrad Lohse. ABLE: blockwise site frequency spectra for inferring complex population histories and recombination. *Genome Biology, Year = 2018, Volume = 19, Month = SEP 25, DOI = 10.1186/s13059-018-1517-y, Article-Number = 145, ISSN = 1474-760X, ORCID-Numbers = Beeravolu Reddy, Champak/0000-0002-0800-1994 Frantz, Laurent/0000-0001-8030-3885, Times-Cited = 3, Unique-ID = ISI:000445752300004,*.

[3]  Annabel C. Beichman, Tanya N. Phung, and Kirk E. Lohmueller. Comparison of Single Genome and Allele Frequency Data Reveals Discordant Demographic Histories. *G3-GENES GENOMES GENETICS*, 7(11):3605–3620, NOV 2017.

[4]  Anders Bergstrom, Shane A. McCarthy, Ruoyun Hui, Mohamed A. Almarri, Qasim Ayub, Petr Danecek, Yuan Chen, Sabine Felkel, Pille Hallast, Jack Kamm, Helene Blanche, Jean-Francois Deleuze, Howard Cann, Swapan Mallick, David Reich, Manjinder S. Sandhu, Pontus Skoglund, Aylwyn Scally, Yali Xue, Richard Durbin, and Chris Tyler-Smith. Insights into human genetic variation and population history from 929 diverse genomes. *SCIENCE*, 367(6484, SI):1339+, MAR 20 2020.

[5]  Sharon R. Browning, Brian L. Browning, Ying Zhou, Serena Tucci, and Joshua M. Akey. Analysis of Human Sequence Data Reveals Two Pulses of Archaic Denisovan Admixture. *CELL*, 173(1):53+, MAR 22 2018.

[6]  Jun Cao, Korbinian Schneeberger, Stephan Ossowski, Torsten Guenther, Sebastian Bender, Joffrey Fitz, Daniel Koenig, Christa Lanz, Oliver Stegle, Christoph Lippert, Xi Wang, Felix Ott, Jonas Mueller, Carlos Alonso-Blanco, Karsten Borg-

42

**684**     wardt, Karl J. Schmid, and Detlef Weigel. Whole-genome sequencing of multiple

**685**     Arabidopsis thaliana populations. *Nature Genetics*, 43(10):956–U60, OCT 2011.

**686**  [7] Dan Chang and Beth Shapiro. Using ancient DNA and coalescent-based methods

**687**     to infer extinction. *Biology Letters*, 12(2), FEB 1 2016.

**688**  [8] Lounes Chikhi, Willy Rodriguez, Simona Grusea, Patricia Santos, Simon Boitard,

**689**     and Olivier Mazet. The IICR (inverse instantaneous coalescence rate) as a sum-

**690**     mary of genomic diversity: insights into demographic inference and model choice.

**691**     *Heredity*, 120(1):13–24, JAN 2018.

**692**  [9] Slew Woh Choo, Mike Rayko, Tze King Tan, Ranjeev Hari, Aleksey Komissarov,

**693**     Wei Yee Wee, Andrey A. Yurchenko, Sergey Kliver, Gaik Tamazian, Agostinho

**694**     Antunes, Richard K. Wilson, Wesley C. Warren, Klaus-Peter Koepfli, Patrick

**695**     Minx, Ksenia Krasheninnikova, Antoinette Kotze, Desire L. Dalton, Elaine Ver-

**696**     maak, Ian C. Paterson, Pavel Dobrynin, Frankie Thomas Sitam, Jeffrine J.

**697**     Rovie-Ryan, Warren E. Johnson, Aini Mohamed Yusoff, Shu-Jin Luo, Kayal Vizi

**698**     Karuppannan, Gang Fang, Deyou Zheng, Mark B. Gerstein, Leonard Lipovich,

**699**     Stephen J. O'Brien, and Guat Jah Wong. Pangolin genomes and the evolution

**700**     of mammalian scales and immunity. *GENOME RESEARCH*, 26(10):1312–1322,

**701**     OCT 2016.

**702** [10] Robert Ekblom, Birte Brechlin, Jens Persson, Linnea Smeds, Malin Johansson,

**703**     Jessica Magnusson, Oystein Flagstad, and Hans Ellegren. Genome sequencing and

**704**     conservation genomics in the Scandinavian wolverine population. *Conservation*

**705**     *Biology*, 32(6):1301–1312, DEC 2018.

**706** [11] Adam D. Ewing. Transposable element detection from whole genome sequence

**707**     data. *MOBILE DNA*, 6, DEC 29 2015.

[12] Andrea Fulgione, Maarten Koornneef, Fabrice Roux, Joachim Hermisson, and Angela M. Hancock. Madeiran Arabidopsis thaliana Reveals Ancient Long-Range Colonization and Clarifies Demography in Eurasia. *Molecular Biology and Evolution*, 35(3):564–574, MAR 2018.

[13] Lucie Gattepaille, Torsten Guenther, and Mattias Jakobsson. Inferring Past Effective Population Size from Distributions of Coalescent Times. *Molecular Biology and Evolution*, 204(3):1191+, NOV 2016.

[14] Brandon S. Gaut, Danelle K. Seymour, Qingpo Liu, and Yongfeng Zhou. Demography and its effects on genomic variation in crop domestication. *Nature Plants*, 4(8):512–520, AUG 2018.

[15] John Hawks. Introgression Makes Waves in Inferred Histories of Effective Population Size. *HUMAN BIOLOGY*, 89(1):67–80, JAN 2017.

[16] Luke B. B. Hecht, Peter C. Thompson, and Benjamin M. Rosenthal. Comparative demography elucidates the longevity of parasitic and symbiotic relationships. *PROCEEDINGS OF THE ROYAL SOCIETY B-BIOLOGICAL SCIENCES*, 285(1888), OCT 10 2018.

[17] Sarah Hendricks, Eric C. Anderson, Tiago Antao, Louis Bernatchez, Brenna R. Forester, Brittany Garner, Brian K. Hand, Paul A. Hohenlohe, Martin Kardos, Ben Koop, Arun Sethuraman, Robin S. Waples, and Gordon Luikart. Recent advances in conservation and population genomics data analysis. *Evolutionary Applications*, 11(8):1197–1211, SEP 2018.

[18] Asger Hobolth and Jens Ledet Jensen. Markovian approximation to the finite

730    loci coalescent with recombination along multiple sequences. *THEORETICAL*

731    *POPULATION BIOLOGY*, 98:48–58, DEC 2014.

732    [19] James E. Johndrow and Julia A. Palacios. Exact limits of inference in coalescent

733    models. *Theoretical Population Biology*, 125:75–93, FEB 2019.

734    [20] Marty Kardos, Anna Qvarnstrom, and Hans Ellegren. Inferring Individual In-

735    breeding and Demographic History from Segments of Identity by Descent in

736    Ficedula Flycatcher Genome Sequences. *GENETICS*, 205(3):1319–1334, MAR

737    2017.

738    [21] Jerome Kelleher, Alison M. Etheridge, and Gilean McVean. Efficient Coalescent

739    Simulation and Genealogical Analysis for Large Sample Sizes. *PLOS COMPU-*

740    *TATIONAL BIOLOGY*, 12(5), MAY 2016.

741    [22] Jerome Kelleher, Yan Wong, Anthony W. Wohns, Chaimaa Fadil, Patrick K.

742    Albers, and Gil McVean. Inferring whole-genome histories in large population

743    datasets (vol 51, pg 1330, 2019). *NATURE GENETICS*, 51(11):1660, NOV 2019.

744    [23] Younhun Kim, Frederic Koehler, Ankur Moitra, Elchanan Mossel, and Govind

745    Ramnarayan. How Many Subpopulations Is Too Many? Exponential Lower

746    Bounds for Inferring Population Histories. *JOURNAL OF COMPUTATIONAL*

747    *BIOLOGY*, 27(4):613–625, APR 1 2020.

748    [24] Robert Kofler. SimulaTE: simulating complex landscapes of transposable ele-

749    ments of populations. *BIOINFORMATICS*, 34(8):1419–1420, APR 15 2018.

750    [25] Sally C. Y. Lau, Nerida G. Wilson, Catarina N. S. Silva, and Jan M. Strugnell.

751    Detecting glacial refugia in the Southern Ocean. *ECOGRAPHY*.

[26] Heng Li and Richard Durbin. Inference of human population history from individual whole-genome sequences. *Nature*, 475(7357):493–U84, JUL 28 2011.

[27] Shengbin Li, Bo Li, Cheng Cheng, Zijun Xiong, Qingbo Liu, Jianghua Lai, Hannah V. Carey, Qiong Zhang, Haibo Zheng, Shuguang Wei, Hongbo Zhang, Liao Chang, Shiping Liu, Shanxin Zhang, Bing Yu, Xiaofan Zeng, Yong Hou, Wenhui Nie, Youmin Guo, Teng Chen, Jiuqiang Han, Jian Wang, Jun Wang, Chen Chen, Jiankang Liu, Peter J. Stambrook, Ming Xu, Guojie Zhang, M. Thomas P. Gilbert, Huanming Yang, Erich D. Jarvis, Jun Yu, and Jianqun Yan. Genomic signatures of near-extinction and rebirth of the crested ibis and other endangered bird species. *GENOME BIOLOGY*, 15(12), 2014.

[28] Michael Lynch, Ryan Gutenkunst, Matthew Ackerman, Ken Spitze, Zhiqiang Ye, Takahiro Maruki, and Zhiyuan Jia. Population Genomics of Daphnia pulex. *Molecular Biology and Evolution*, 206(1):315–332, MAY 2017.

[29] Michael Lynch, Bernhard Haubold, Peter Pfaffelhuber, and Takahiro Maruki. Inference of Historical Population-Size Changes with Allele-Frequency Data. *G3-GENES GENOMES GENETICS*, 10(1):211–223, JAN 2020.

[30] Anna-Sapfo Malaspinas, Michael C. Westaway, Craig Muller, Vitor C. Sousa, Oscar Lao, Isabel Alves, Anders Bergstrom, Georgios Athanasiadis, Jade Y. Cheng, Jacob E. Crawford, Tim H. Heupink, Enrico Macholdt, Stephan Peischl, Simon Rasmussen, Stephan Schiffels, Sankar Subramanian, Joanne L. Wright, Anders Albrechtsen, Chiara Barbieri, Isabelle Dupanloup, Anders Eriksson, Ashot Margaryan, Ida Moltke, Irina Pugach, Thorfinn S. Korneliussen, Ivan P. Levkivskyi, J. Vctor Moreno-Mayar, Shengyu Ni, Fernando Racimo, Martin Sikora, Yali Xue, Farhang A. Aghakhanian, Nicolas Brucato, Soren Brunak,

46

Paula F. Campos, Warren Clark, Sturla Ellingvag, Gudjugudju Fourmile, Pascale Gerbault, Darren Injie, George Koki, Matthew Leavesley, Betty Logan, Aubrey Lynch, Elizabeth A. Matisoo-Smith, Peter J. McAllister, Alexander J. Mentzer, Mait Metspalu, Andrea B. Migliano, Les Murgha, Maude E. Phipps, William Pomat, Doc Reynolds, Francois-Xavier Ricaut, Peter Siba, Mark G. Thomas, Thomas Wales, Colleen Ma'run Wall, Stephen J. Oppenheimer, Chris Tyler-Smith, Richard Durbin, Joe Dortch, Andrea Manica, Mikkel H. Schierup, Robert A. Foley, Marta Mirazon Lahr, Claire Bowern, Jeffrey D. Wall, Thomas Mailund, Mark Stoneking, Rasmus Nielsen, Manjinder S. Sandhu, Laurent Excoffier, David M. Lambert, and Eske Willerslev. A genomic history of Aboriginal Australia. *NATURE*, 538(7624):207+, OCT 13 2016.

[31] P Marjoram and JD Wall. Fast "coalescent" simulation. *BMC Genetics*, 7, MAR 15 2006.

[32] Niklas Mather, Samuel M. Traves, and Simon Y. W. Ho. A practical introduction to sequentially Markovian coalescent methods for estimating demographic history from genomic data. *ECOLOGY AND EVOLUTION*, 10(1):579–589, JAN 2020.

[33] Maja P. Mattle-Greminger, Tugce Bilgin Sonay, Alexander Nater, Marc Pybus, Tariq Desai, Guillem de Valles, Ferran Casals, Aylwyn Scally, Jaume Bertranpetit, Tomas Marques-Bonet, Carel P. van Schaik, Maria Anisimova, and Michael Kruetzen. Genomes reveal marked differences in the adaptive evolution between orangutan species. *Genome Biology*, 19, NOV 15 2018.

[34] O. Mazet, W. Rodriguez, S. Grusea, S. Boitard, and L. Chikhi. On the importance of being structured: instantaneous coalescence rates and human evolution-lessons for ancestral population size inference? *Heredity*, 116(4):362–371, APR 2016.

47

[35] GAT McVean and NJ Cardin. Approximating the coalescent with recombination. *Philosophical Transactions of the Royal Society B-Biological Sciences*, 360(1459):1387–1393, JUL 29 2005.

[36] Sajad Mirzaei and Yufeng Wu. RENT plus : an improved method for inferring local genealogical trees from haplotypes with recombination. *BIOINFORMATICS*, 33(7):1021–1030, APR 1 2017.

[37] Krystyna Nadachowska-Brzyska, Reto Burri, Linnea Smeds, and Hans Ellegren. PSMC analysis of effective population sizes in molecular ecology and its application to black-and-white Ficedula flycatchers. *Molecular Ecology*, 25(5):1058–1072, MAR 2016.

[38] Shigeki Nakagome, Richard R. Hudson, and Anna Di Rienzo. Inferring the model and onset of natural selection under varying population size from the site frequency spectrum and haplotype structure. *PROCEEDINGS OF THE ROYAL SOCIETY B-BIOLOGICAL SCIENCES*, 286(1896), FEB 6 2019.

[39] Michael G. Nelson, Raquel S. Linheiro, and Casey M. Bergman. McClintock: An Integrated Pipeline for Detecting Transposable Element Insertions in Whole-Genome Shotgun Sequencing Data. *G3-GENES GENOMES GENETICS*, 7(8):2763–2778, AUG 2017.

[40] Kevin P. Oh, Cameron L. Aldridge, Jennifer S. Forbey, Carolyn Y. Dadabay, and Sara J. Oyler-McCance. Conservation Genomics in the Sagebrush Sea: Population Divergence, Demographic History, and Local Adaptation in Sage-Grouse (Centrocercus spp.). *GENOME BIOLOGY AND EVOLUTION*, 11(7):2023–2034, JUL 2019.

48

[41] Julia A. Palacios, John Wakeley, and Sohini Ramachandran. Bayesian Nonparametric Inference of Population Size Changes from Sequential Genealogies. *Genetics*, 201(1):281+, SEP 2015.

[42] Pier Francesco Palamara, Jonathan Terhorst, Yun S. Song, and Alkes L. Price. High-throughput inference of pairwise coalescence times identifies signals of selection and enriched disease heritability. *NATURE GENETICS*, 50(9):1311+, SEP 2018.

[43] Eleftheria Palkopoulou, Mark Lipson, Swapan Mallick, Svend Nielsen, Nadin Rohland, Sina Baleka, Emil Karpinski, Atma M. Ivancevici, Thu-Hien To, Daniel Kortschak, Joy M. Raison, Zhipeng Qu, Tat-Jun Chin, Kurt W. Alt, Stefan Claesson, Love Dalen, Ross D. E. MacPhee, Harald Meller, Alfred L. Rocar, Oliver A. Ryder, David Heiman, Sarah Young, Matthew Breen, Christina Williams, Bronwen L. Aken, Magali Ruffier, Elinor Karlsson, Jeremy Johnson, Federica Di Palma, Jessica Alfoldi, David L. Adelsoni, Thomas Mailund, Kasper Munch, Kerstin Lindblad-Toh, Michael Hofreiter, Hendrik Poinar, and David Reich. A comprehensive genomic history of extinct and living elephants. *Proceedings of the National Academy of Sciences of the United States of America*, 115(11):E2566–E2574, MAR 13 2018.

[44] Eleftheria Palkopoulou, Swapan Mallick, Pontus Skoglund, Jacob Enk, Nadin Rohland, Heng Li, Ayca Omrak, Sergey Vartanyan, Hendrik Poinar, Anders Gotherstrom, David Reich, and Love Dalen. Complete Genomes Reveal Signatures of Demographic and Genetic Declines in the Woolly Mammoth. *Current Biology*, 25(10):1395–1400, MAY 18 2015.

[45] Austin H. Patton, Mark J. Margres, Amanda R. Stahlke, Sarah Hendricks, Kevin

49

Lewallen, Rodrigo K. Hamede, Manuel Ruiz-Aravena, Oliver Ryder, Hamish Mc-Callum, I, Menna E. Jones, Paul A. Hohenlohe, and Andrew Storfer. Contemporary Demographic Reconstruction Methods Are Robust to Genome Assembly Quality: A Case Study in Tasmanian Devils. *MOLECULAR BIOLOGY AND EVOLUTION*, 36(12):2906–2921, DEC 2019.

[46] S. P. Pfeifer. From next-generation resequencing reads to a high-quality variant data set. *HEREDITY*, 118(2):111–124, FEB 2017.

[47] Roy N. Platt, II, Laura Blanco-Berdugo, and David A. Ray. Accurate Transposable Element Annotation Is Vital When Analyzing New Genome Assemblies. *GENOME BIOLOGY AND EVOLUTION*, 8(2):403–410, FEB 2016.

[48] Javier Prado-Martinez, Peter H. Sudmant, Jeffrey M. Kidd, Heng Li, Joanna L. Kelley, Belen Lorente-Galdos, Krishna R. Veeramah, August E. Woerner, Timothy D. O'Connor, Gabriel Santpere, Alexander Cagan, Christoph Theunert, Ferran Casals, Hafid Laayouni, Kasper Munch, Asger Hobolth, Anders E. Halager, Maika Malig, Jessica Hernandez-Rodriguez, Irene Hernando-Herraez, Kay Pruefer, Marc Pybus, Laurel Johnstone, Michael Lachmann, Can Alkan, Dorina Twigg, Natalia Petit, Carl Baker, Fereydoun Hormozdiari, Marcos Fernandez-Callejo, Marc Dabad, Michael L. Wilson, Laurie Stevison, Cristina Camprubi, Tiago Carvalho, Aurora Ruiz-Herrera, Laura Vives, Marta Mele, Teresa Abello, Ivanela Kondova, Ronald E. Bontrop, Anne Pusey, Felix Lankester, John A. Kiyang, Richard A. Bergl, Elizabeth Lonsdorf, Simon Myers, Mario Ventura, Pascal Gagneux, David Comas, Hans Siegismund, Julie Blanc, Lidia Agueda-Calpena, Marta Gut, Lucinda Fulton, Sarah A. Tishkoff, James C. Mullikin, Richard K. Wilson, Ivo G. Gut, Mary Katherine Gonder, Oliver A. Ryder, Beatrice H. Hahn,

50

871    Arcadi Navarro, Joshua M. Akey, Jaume Bertranpetit, David Reich, Thomas

872    Mailund, Mikkel H. Schierup, Christina Hvilsom, Aida M. Andres, Jeffrey D.

873    Wall, Carlos D. Bustamante, Michael F. Hammer, Evan E. Eichler, and Tomas

874    Marques-Bonet. Great ape genetic diversity and population history. *NATURE*,

875    499(7459):471–475, JUL 25 2013.

876  [49] Willy Rodriguez, Olivier Mazet, Simona Grusea, Armando Arredondo, Josue M.

877    Corujo, Simon Boitard, and Lounes Chikhi. The IICR and the non-stationary

878    structured coalescent: towards demographic inference with arbitrary changes in

879    population structure. *Heredity*, 121(6):663–678, DEC 2018.

880  [50] Stephan Schiffels and Richard Durbin. Inferring human population size and sep-

881    aration history from multiple genome sequences. *Nature Genetics*, 46(8):919–925,

882    AUG 2014.

883  [51] Joshua G. Schraiber and Joshua M. Akey. Methods and models for unravelling

884    human evolutionary history. *NATURE REVIEWS GENETICS*, 16(12):727–740,

885    DEC 2015.

886  [52] Daniel R. Schrider, Alexander G. Shanku, and Andrew D. Kern. Effects of Linked

887    Selective Sweeps on Demographic Inference and Model Selection. *GENETICS*,

888    204(3):1207+, NOV 2016.

889  [53] Thibaut Paul Patrick Sellinger, Diala Abu Awad, Markus Moest, and Aurelien

890    Tellier. Inference of past demography, dormancy and self-fertilization rates from

891    whole genome sequence data. *PLOS GENETICS*, 16(4), APR 2020.

892  [54] Sara Sheehan, Kelley Harris, and Yun S. Song. Estimating Variable Effective

893    Population Sizes from Multiple Genomes: A Sequentially Markov Conditional

51

894    Sampling Distribution Approach. *Molecular Biology and Evolution*, 194(3):647+,

895    JUL 2013.

896    [55] Sara Sheehan and Yun S. Song. Deep Learning for Population Genetic Inference.

897    *PLOS Computational Biology*, 12(3), MAR 2016.

898    [56] Montgomery Slatkin. Statistical methods for analyzing ancient DNA from ho-

899    minins. *CURRENT OPINION IN GENETICS & DEVELOPMENT*, 41:72–76,

900    DEC 2016.

901    [57] Chris C. R. Smith and Samuel M. Flaxman. Leveraging whole genome sequenc-

902    ing data for demographic inference with approximate Bayesian computation.

903    *MOLECULAR ECOLOGY RESOURCES*, 20(1):125–139, JAN 2020.

904    [58] Leo Speidel, Marie Forest, Sinan Shi, and Simon R. Myers. A method for genome-

905    wide genealogy estimation for thousands of samples. *NATURE GENETICS*,

906    51(9):1321+, SEP 2019.

907    [59] Jeffrey P. Spence, Matthias Steinrucken, Jonathan Terhorst, and Yun S. Song.

908    Inference of population history using coalescent HMMs: review and outlook. *Cur-*

909    *rent Opinion in Genetics & Development*, 53:70–76, DEC 2018.

910    [60] Paul R. Staab, Sha Zhu, Dirk Metzler, and Gerton Lunter. scrm: efficiently

911    simulating long sequences using the approximated coalescent with recombination.

912    *Bioinformatics*, 31(10):1680–1682, MAY 15 2015.

913    [61] Remco Stam, Tetyana Nosenko, Anja C. Hoerger, Wolfgang Stephan, Michael

914    Seidel, Jose M. M. Kuhn, Georg Haberer, and Aurelien Tellier. The de Novo

915    Reference Genome and Transcriptome Assemblies of the Wild Tomato Species

52

Solanum chilense Highlights Birth and Death of NLR Genes Between Tomato Species. *G3-GENES GENOMES GENETICS*, 9(12):3933–3941, DEC 2019.

[62] Matthias Steinrucken, Jack Kamm, Jeffrey P. Spence, and Yun S. Song. Inference of complex population histories using whole-genome sequences from multiple populations. *PROCEEDINGS OF THE NATIONAL ACADEMY OF SCIENCES OF THE UNITED STATES OF AMERICA*, 116(34):17115–17120, AUG 20 2019.

[63] Jonathan Terhorst, John A. Kamm, and Yun S. Song. Robust and scalable inference of population history froth hundreds of unphased whole genomes. *Nature Genetics*, 49(2):303–309, FEB 2017.

[64] Jonathan Terhorst and Yun S. Song. Fundamental limits on the accuracy of demographic inference based on the sample frequency spectrum. *Proceedings of the National Academy of Sciences of the United States of America*, 112(25):7677–7682, JUN 23 2015.

[65] Berit Lindum Waltoft and Asger Hobolth. Non-parametric estimation of population size changes from the site frequency spectrum. *Statistical Applications in Genetics and Molecular Biology*, 17(3), JUN 2018.

[66] Ke Wang, Iain Mathieson, Jared O'Connell, and Stephan Schiffels. Tracking human population structure through time from whole genome sequences. *PLOS GENETICS*, 16(3), MAR 2020.

[67] Pengcheng Wang, Hongyan Yao, Kadeem J. Gilbert, Qi Lu, Yu Hao, Zhengwang Zhang, and Nan Wang. Glaciation-based isolation contributed to speciation in a Palearctic alpine biodiversity hotspot: Evidence from endemic species. *Molecular Phylogenetics and Evolution*, 129:315–324, DEC 2018.

53

[68] Rachel C. Williams, Marina B. Blanco, Jelmer W. Poelstra, Kelsie E. Hunnicutt, Aaron A. Comeault, and Anne D. Yoder. Conservation genomic analysis reveals ancient introgression and declining levels of genetic diversity in Madagascar's hibernating dwarf lemurs. *HEREDITY*, 124(1):236–251, JAN 2020.

[69] C Wiuf and J Hein. Recombination as a point process along sequences. *Theoretical Population Biology*, 55(3):248–259, JUN 1999.

[70] Chee-Wei Yew, Dongsheng Lu, Lian Deng, Lai-Ping Wong, Rick Twee-Hee Ong, Yan Lu, Xiaoji Wang, Yushimah Yunus, Farhang Aghakhanian, Siti Shuhada Mokhtar, Mohammad Zahirul Hoque, Christopher Lok-Yung Voo, Thuhairah Abdul Rahman, Jong Bhak, Maude E. Phipps, Shuhua Xu, Yik-Ying Teo, Subbiah Vijay Kumar, and Boon-Peng Hoh. Genomic structure of the native inhabitants of Peninsular Malaysia and North Borneo suggests complex human population history in Southeast Asia. *Human Genetics*, 137(2):161–173, FEB 2018.